# CS 744: PIPEDREAM

Shivaram Venkataraman
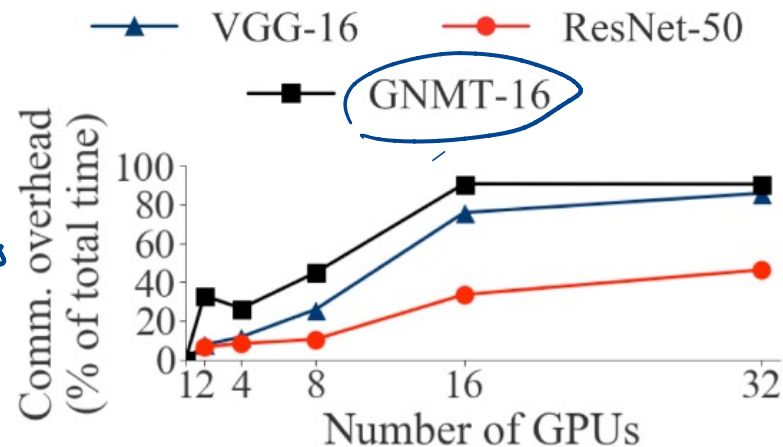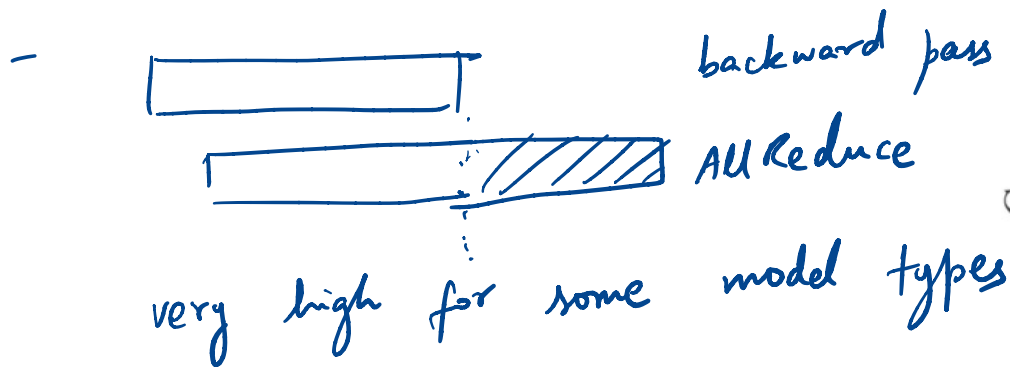
Fall 2022

# ADMINISTRIVIA

- Assignment 2 is due Wednesday AM! → *Please post on Piazza*

- Course project preference sheet: Out today!
    - Propose your own?
    - Or rank 1 through 5 of some project ideas we have

    – *Group!*

# LIMITATIONS OF DATA PARALLEL

- Overhead increases as num GPUs increases

-

backward pass

AllReduce

very high for some model types



8xV100s with NVLink (AWS)
PyTorch + NCCL 2.4

- Overhead is lower <= 8 GPUs, higher after that

1 to 8 GPUs

machine

8 GPUs

machine

"fraction of training time spent in communication stalls"

# MODEL PARALLEL TRAINING

5 layers in model

low utilization

stage 1

**Worker 1  Worker 2  Worker 3      Worker 4**

input

Input stage

Output stage

stage 4

forward pass layer 1

fwd pass layer 2



| | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Worker 1 | 1 | | | | | | | | | 1 | 1 | 2 | | |
| Worker 2 | | 1 | | | | | | 1 | 1 | | | | 2 | |
| Worker 3 | | | 1 | | | | 1 | 1 | | | | | | 2 |
| Worker 4 | | | | 1 | 1 | 1 | | | | | | | | 2 |

Time

| Forward Pass | Backward Pass | Idle |
|---|---|---|

backward pass for layer 4,5

gradients here

① Comm is constant as num workers increases. limited connection

② Memory req. are lower & con use heterogeneous machines.

# PIPELINE PARALLEL
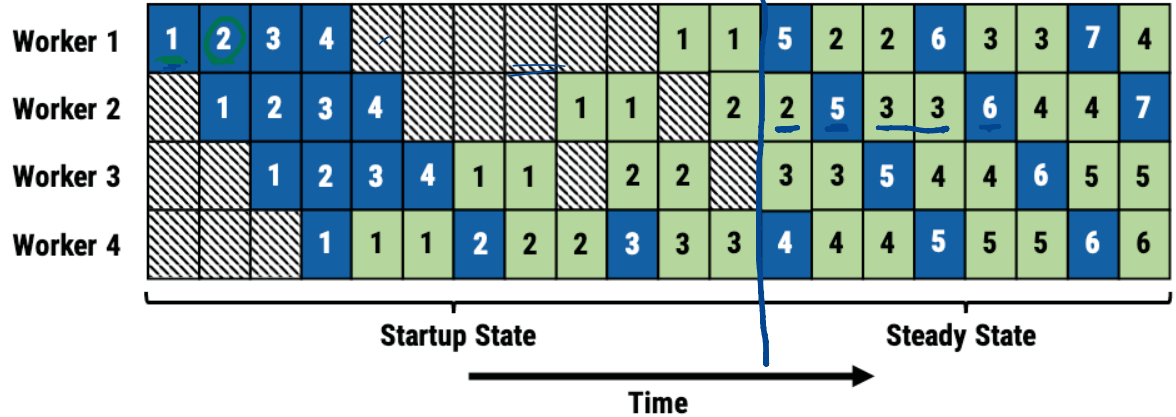
Key Idea:
Instead of 1 input batch
feed in "k" batches

- Handle scenarios
where some stages
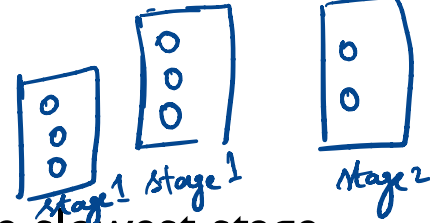are slow.

$k = 4$ here

steady state
util is high



| | Worker 1 | 1 | 2 | 3 | 4 | | | | | | | 1 | 1 | 5 | 2 | 2 | 6 | 3 | 3 | 7 | 4 |

Advantages?

- fixed comm. overhead
- high utilization
-

- Partitioning
- Scheduling of FWD, BCKWD
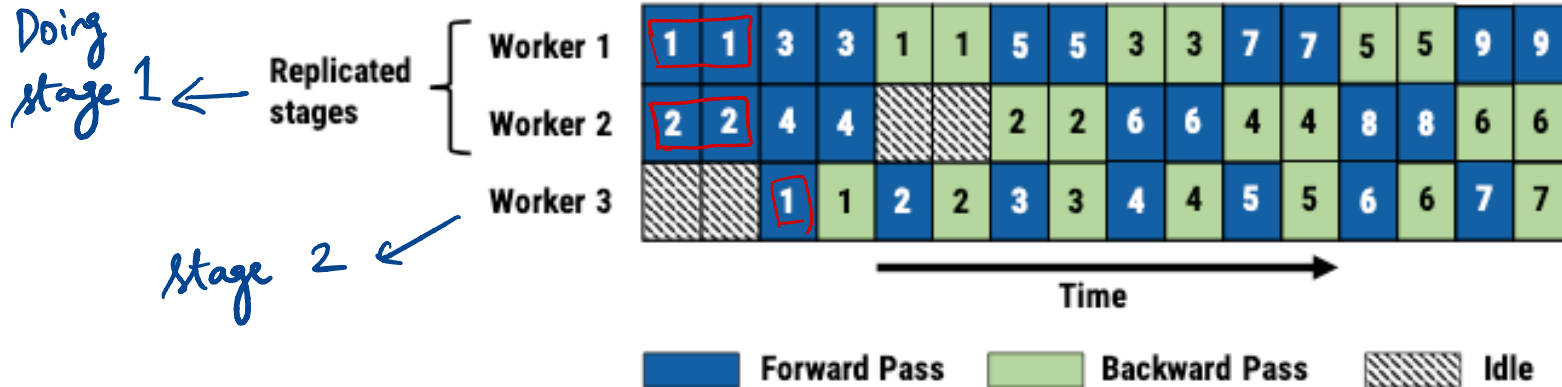- Learning

# CHALLENGE 1: WORK PARTITIONING

*similar time?*

Goal: Balanced stages in the pipeline. Why?

Steady state throughput is the throughput of the slowest stage

*Hybrid parallel*

Stages can be replicated! Ex: Two stage pipeline, but first stage is replicated

*Doing stage 1* ←

*Stage 2* ←



| | | Time → | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Worker 1 (Replicated stages) | 1 | 1 | 3 | 3 | 1 | 1 | 5 | 5 | 3 | 3 | 7 | 7 | 5 | 5 | 9 | 9 |
| Worker 2 (Replicated stages) | 2 | 2 | 4 | 4 | | | 2 | 2 | 6 | 6 | 4 | 4 | 8 | 8 | 6 | 6 |
| Worker 3 | | | 1 | 1 | 2 | 2 | 3 | 3 | 4 | 4 | 5 | 5 | 6 | 6 | 7 | 7 |

Forward Pass    Backward Pass    Idle

# WORK PARITIONING

Profiler: computation time for forward, backward for each layer

size of output activations, gradients (network transfer)

size of parameters (memory)

Dynamic programming algorithm

Intuition: Find optimal partitions within a server,

Then find best split across servers using that

Given
- model

- cluster
  ↳ GPUs, network

Static   plannings

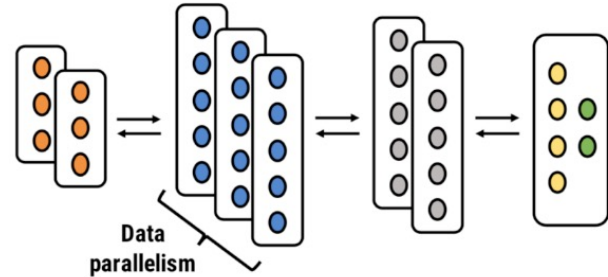# CHALLENGE 2: WORK SCHEDULING

Traditional data parallel

    forward iter(i)     *Iteration 0*

    backward iter(i)

    forward iter(i+1)     *Iteration 1*

    …

Pipeline parallel: Worker can

    Forward pass to push to downstream
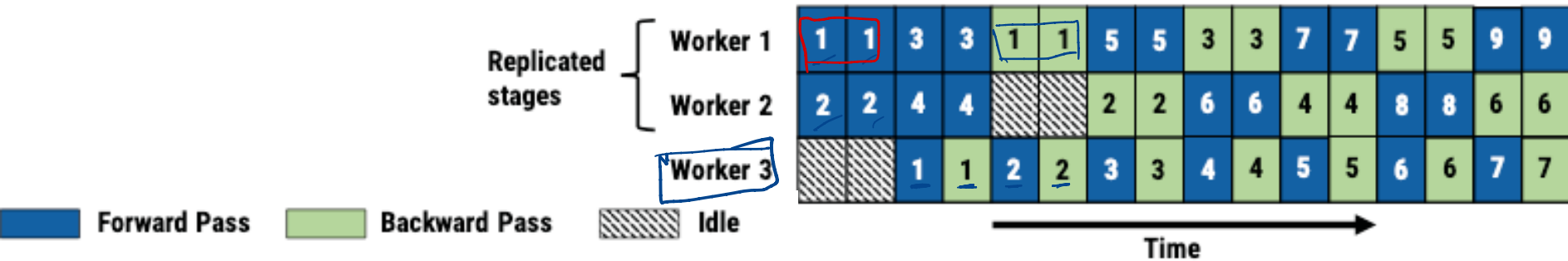
    Backward pass to push to upstream



Data parallelism

*Either run fwd pass next batch or backward pass for prior batch*

# CHALLENGE 2: WORK SCHEDULING

Num active batches ~= num_workers / num_replicas_input

Schedule one-forward-one-backward (1F1B) – Worker 3 → *makes sure pipeline is making progress*

Round-robin for replicated stages → Worker 2

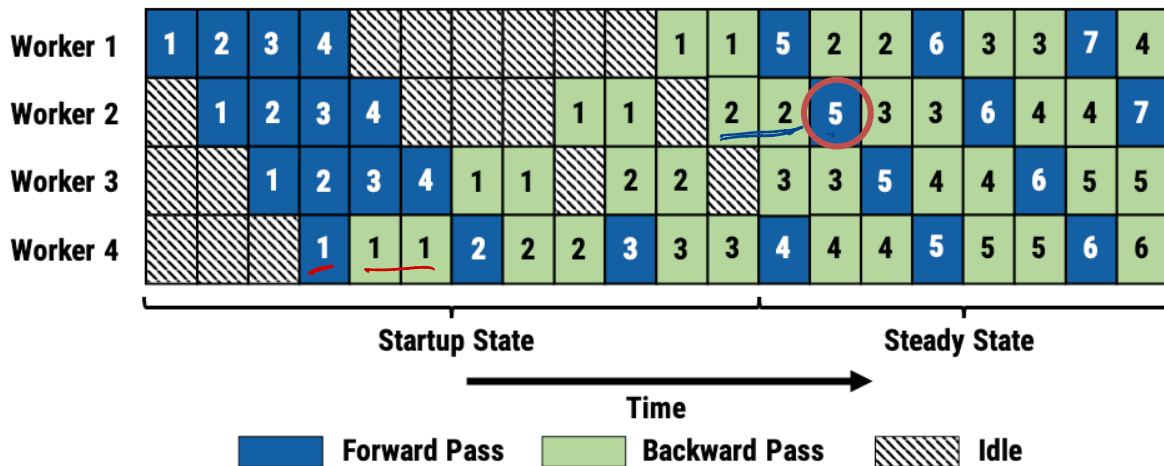    same worker for fwd, backward

# CHALLENGE 3: EFFECTIVE LEARNING

Naïve pipelining

Different model versions forward and backward

Batch no. 5
↳ fwd pass uses model r2.
↳ bck pass used model v4.



| | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Worker 1 | 1 | 2 | 3 | 4 | | | | | | 1 | 1 | 5 | 2 | 2 | 6 | 3 | 3 | 7 | 4 |
| Worker 2 | | 1 | 2 | 3 | 4 | | | 1 | 1 | | 2 | 2 | 5 | 3 | 3 | 6 | 4 | 4 | 7 |
| Worker 3 | | | 1 | 2 | 3 | 4 | 1 | 1 | | 2 | 2 | | 3 | 3 | 5 | 4 | 4 | 6 | 5 | 5 |
| Worker 4 | | | | 1 | 1 | 1 | 2 | 2 | 2 | 3 | 3 | 3 | 4 | 4 | 4 | 5 | 5 | 5 | 6 | 6 |

Startup State       Steady State

Time

Forward Pass    Backward Pass    Idle

Keep model v2 stashed

batch 5 ≡ v2

# CHALLENGE 3: EFFECTIVE LEARNING

batch
no. 5 $= \begin{bmatrix} 32 \\ images \end{bmatrix} \rightarrow m_{V1} \rightarrow$ Activations $\rightarrow m_{V_2} \rightarrow$ Acti

$\underbrace{\qquad}_{W1}$    $\underbrace{\qquad}_{V_2}$
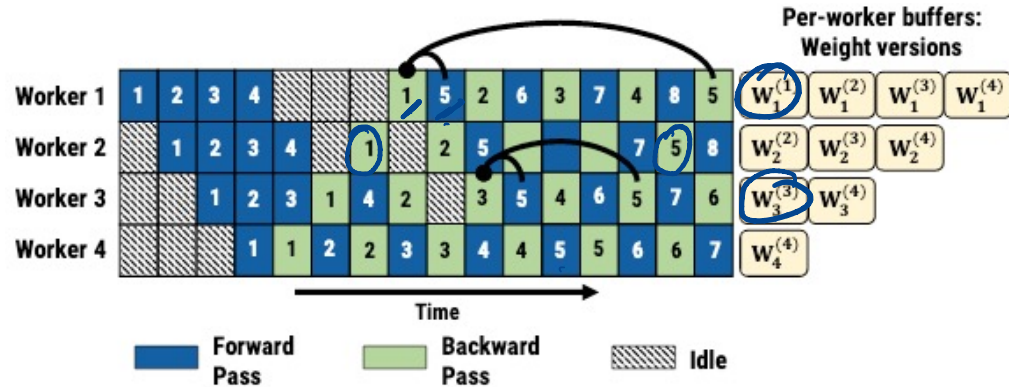
Weight stashing

    Maintain multiple versions of the weights

    One per active mini-batch

Use latest version for forward pass.

Retrieve for backward

No guarantees across stages!

$\hookrightarrow$ diff model versions used in diff stages



Per-worker buffers:
Weight versions

| Worker 1 | Worker 2 | Worker 3 | Worker 4 |

Time

Forward Pass    Backward Pass    Idle

# STALENESS, MEMORY OVERHEAD

How to avoid staleness:

Vertical sync $\longrightarrow$ in the first stage, record model version used propagate that across workers

Memory overhead

Similar to data parallel?

$\hookrightarrow$ lose this benefit with vertical sync.

# SUMMARY

Pipeline parallelism: Combine inter-batch and intra-batch

Partitioning: Replication, dynamic programming

Scheduling: 1F1B

Weight management: Stashing, vertical sync

# DISCUSSION

https://forms.gle/5cf16BWn6Dziey6e6
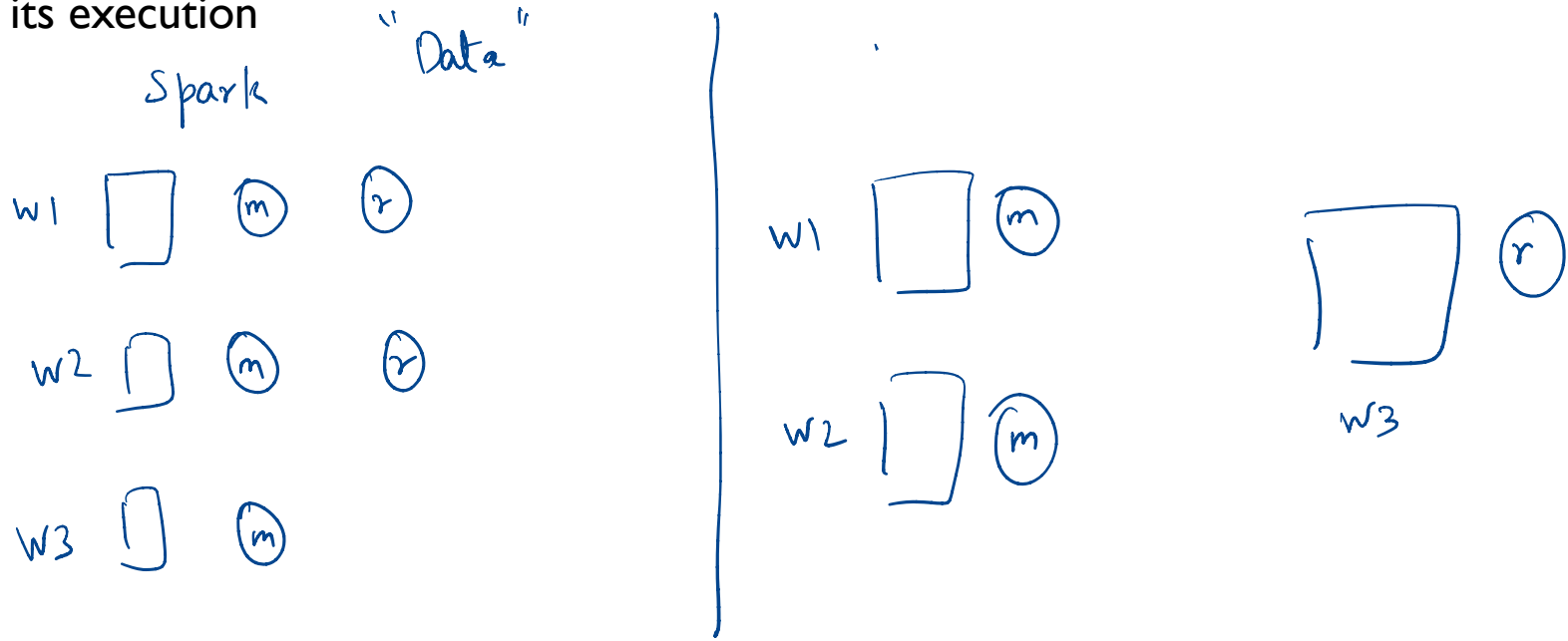
Data parallel - not always bad

4x4 → more overhead for data parallel

List two takeaways from the following table

- speedup higher
VGG-16

| Model Name | Model Size | GPUs (#Servers x #GPUs/Server) | PipeDream Config | Speedup over DataParallel (Epoch Time) |
|---|---|---|---|---|
| Resnet-50 | 97MB | 4x4 2x8 | 16 16 | 1× 1x |
| VGG-16 | 528MB | 4x4 2x8 | 15-1 15-1 | 5.28x 2.98x |
| GNMT-8 | 1.1GB | 3x4 2x8 | Straight 16 | 2.95x 1x |

Why ??

What are some other workload scenarios (e.g. things we discussed for MapReduce or Spark) that could use similar ideas of pipelined parallelism? Develop such one example and its execution

Spark

"Data"

W1 [ ] (m) (r)

W2 [ ] (m) (r)

W3 [ ] (m)

W1 [ ] (m)

W2 [ ] (m)

[ ] (r)

W3

# NEXT STEPS

Next class: Parameter Server

Assignment 2 is due soon!

Course project preference form out today!